# A State-of-the-Art
# Local Training Methods in Federated Learning

**Michal Staňo, Ladislav Hluchý**

INSTITUTE OF INFORMATICS
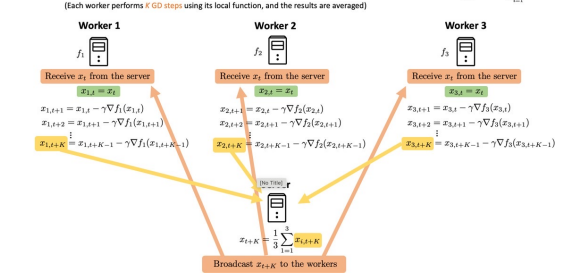SLOVAK ACADEMY OF SCIENCES, BRATISLAVA SLOVAKIA

# Outline of the Talk

1. What is Federated Learning?

2. What is Local Training

3. Brief History of Local Training

4. What does Local Training do?

# Part 1
# What is Federated Learning?

# The First Federated Learning App: Next-Word Prediction?

Federated Learning is a collaborative machine learning from private data stored across a (large) number of clients/devices (e.g., hospitals, phones, banks)

# Part 2
# What is Local Training?

# Local Training

A. Gradient Descent

B. Distributed Gradient Descent

C. Distributed Local Gradient Descent

# Gradient
# Descent

$$\min_{x \in R^d} f(x)$$

$$x_{t+1} = x_t - \gamma \, \nabla f(x_t)$$



$x_t$

$x_{t+1}$

$d=2$

# Distributed Gradient Descent

$$\min_{x \in R^d} f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$

# Distributed Gradient Descent

$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\mathrm{def}}{=} \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$

$n$ — # devices / machines

# model parameters / features

Loss on local data $D_i$ stored on device $i$

$$f_i(x) = \mathbb{E}_{\leftarrow\cdots\leftarrow D_i} f_{i,\leftarrow}(x)$$

The datasets $D_1, \ldots, D_n$ can be arbitrarily heterogeneous

# Distributed Gradient Descent

**Site 1 - Computer**

$f_1$

Receive $x_t$ from the Server

$x_{1,t} = x_t$

$x_{1,t+1} = x_{1,t} - \gamma \nabla f_1(x_{1,t})$

**Site 2 – Mobile Device**

$f_2$

Receive $x_t$ from the Server

$x_{2,t} = x_t$

$x_{2,t+1} = x_{2,t} - \gamma \nabla f_2(x_{2,t})$

**Site 3 - Hospital**

$f_3$

Receive $x_t$ from the Server

$x_{3,t} = x_t$

$x_{3,t+1} = x_{3,t} - \gamma \nabla f_3(x_{3,t})$

**Server**

$$x_{t+1} = \frac{1}{3} \sum_{1=1}^{3} x_{i,t+1} =$$

Broadcast $x_{t+1}$ to the Sites

# Distributed **Local** Gradient Descent

**Site 1,2,3 – Computer, Mobile Device, Hospital**
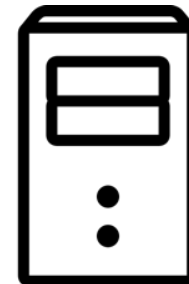
$f_1$ 

Receive $x_t$ from the Server

$x_{1,t} = x_t$

$x_{1,t+1} = x_{1,t} - \gamma \nabla f_1(x_{1,t})$

$x_{1,t+K} = x_{1,t+K-1} - \gamma \nabla f_1(x_{1,t+K-1})$

**Server**
Central Orchestrator

$x_{t+K} = \dfrac{1}{3} \displaystyle\sum_{1-1}^{3} x_{i,t+K}$

Broadcast $x_{t+K}$ to the Sites

# Part 3
# Brief History of Local Training

From Gradient Descent to Local Gradient Descent

**First general Theory for Local GD**
First Analysis of Local GD on Heterogeneous Data
Khaled, Mishchenko, Richtárik

2020

**Federated Averaging - Local GD**
Communication-efficient Learning of Deep Networks from Decentralized Data
H. B. McMahan et all

2017

**Local Gradient Descent Proposed**
Parallel Gradient Distribution in Unconslrained Optimization
O. L. Mangasarian

1995

**Gradient Descent**
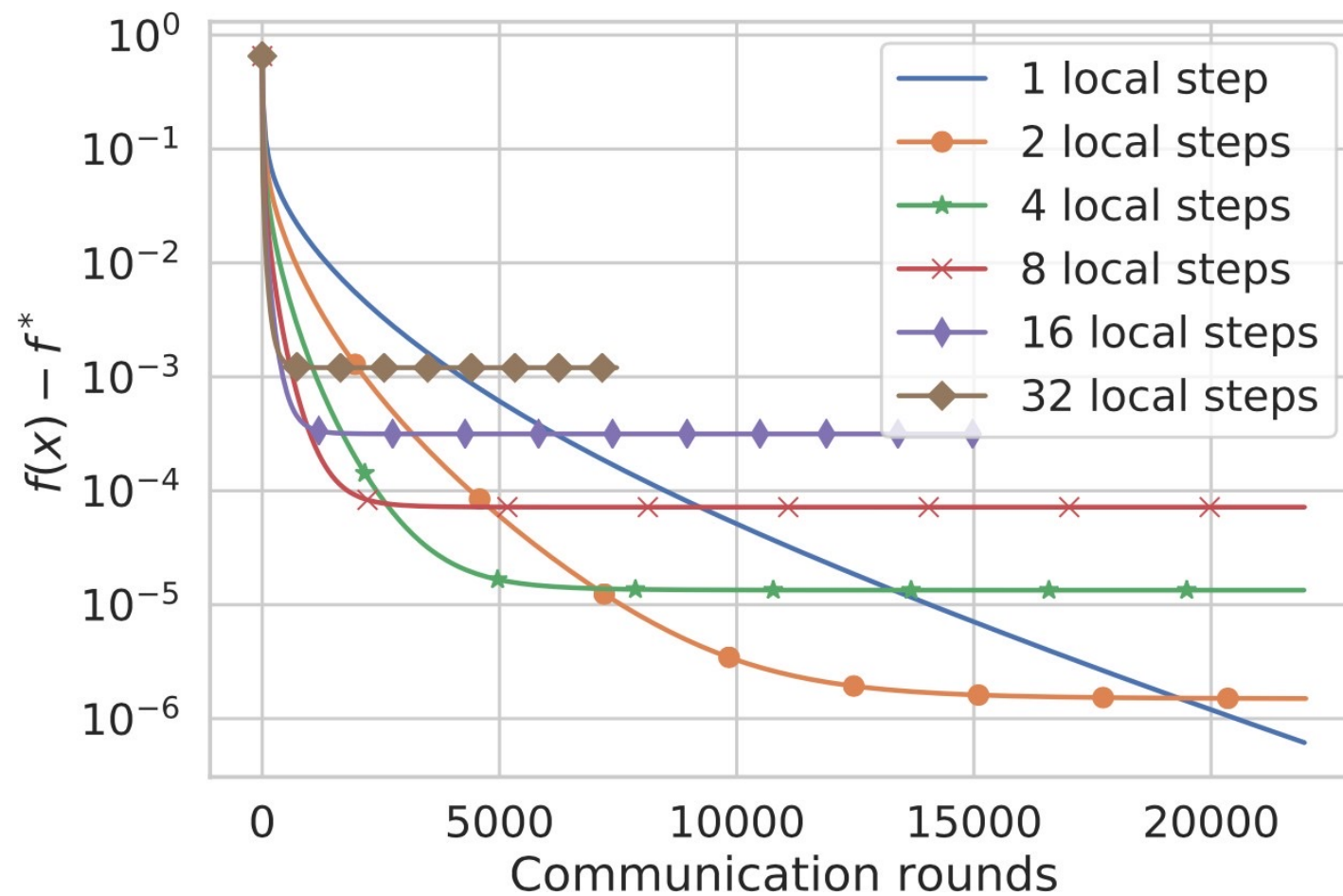Compte Rendu á l' Académie des Sciences
L. A. Cauchy

1847

# Part 4
# What does Local Training do?

# Local Training

# Thank You